

Ignoring the Non-ignorables? Missingness and Missing Positions

König, Thomas; Finke, Daniel; Daimer, Stephanie

Postprint / Postprint

Zeitschriftenartikel / journal article

Zur Verfügung gestellt in Kooperation mit / provided in cooperation with:

www.peerproject.eu

Empfohlene Zitierung / Suggested Citation:

König, T., Finke, D., & Daimer, S. (2005). Ignoring the Non-ignorables? Missingness and Missing Positions. *European Union Politics*, 6(3), 269-290. <https://doi.org/10.1177/1465116505054833>

Nutzungsbedingungen:

Dieser Text wird unter dem "PEER Licence Agreement zur Verfügung" gestellt. Nähere Auskünfte zum PEER-Projekt finden Sie hier: <http://www.peerproject.eu> Gewährt wird ein nicht exklusives, nicht übertragbares, persönliches und beschränktes Recht auf Nutzung dieses Dokuments. Dieses Dokument ist ausschließlich für den persönlichen, nicht-kommerziellen Gebrauch bestimmt. Auf sämtlichen Kopien dieses Dokuments müssen alle Urheberrechtshinweise und sonstigen Hinweise auf gesetzlichen Schutz beibehalten werden. Sie dürfen dieses Dokument nicht in irgendeiner Weise abändern, noch dürfen Sie dieses Dokument für öffentliche oder kommerzielle Zwecke vervielfältigen, öffentlich ausstellen, aufführen, vertreiben oder anderweitig nutzen.

Mit der Verwendung dieses Dokuments erkennen Sie die Nutzungsbedingungen an.

gesis
Leibniz-Institut
für Sozialwissenschaften

Terms of use:

This document is made available under the "PEER Licence Agreement". For more Information regarding the PEER-project see: <http://www.peerproject.eu> This document is solely intended for your personal, non-commercial use. All of the copies of this documents must retain all copyright information and other information regarding legal protection. You are not allowed to alter this document in any way, to copy it for public or commercial purposes, to exhibit the document in public, to perform, distribute or otherwise use the document in public.

By using this particular document, you accept the above-stated conditions of use.

Mitglied der

Leibniz-Gemeinschaft



European Union Politics

DOI: 10.1177/1465116505054833

Volume 6 (3): 269–290

Copyright© 2005

SAGE Publications

London, Thousand Oaks CA,

New Delhi

Ignoring the Non-ignorables?

Missingness and Missing Positions

◆ **Thomas König**

German University of Administrative Sciences, Germany

◆ **Daniel Finke**

Research Institute for Public Administration, Germany

◆ **Stephanie Daimer**

Research Institute for Public Administration, Germany

ABSTRACT

Missing or incomplete data on actors' positions can cause significant problems in political analysis. Research on missing values suggests the use of multiple imputation methods rather than case deletion, but few studies have yet considered the non-ignorable problem – positions that are hidden for strategic purposes. We examine this problem and discuss the advantages and drawbacks of (i) multiple imputation as implemented in AMELIA; (ii) a computationally easy but, in the context of spatial modelling, straightforward measure of indifference and (iii) a conditional averaging algorithm, LDM, which seeks to reasonably fix actors' positions in the policy space pre- and post-imputation. The analysis suggests that actors biased by the status quo strategically hide their more supportive positions. Although none of the existing methods – which produce quite different results – is perfectly suited for imputing hidden positions, LDM has the highest hit rate for the conjectured more supportive position.

KEY WORDS

- EU Constitution
- imputation
- missing values
- strategic positions

Missing positions: An ignorable problem in political analysis?

Missing or incomplete data on actors' bargaining or policy positions can cause significant problems in political analysis. Although empirical political studies can hardly avoid missing positions – whether they gather this information from documents or from expert or mass surveys – we have little knowledge about how to cope with this missing value problem. This deficit is surprising because our theories and tools of analysis are not designed for disregarding missing positions. Ignoring this problem risks biasing our findings: explorative analyses of the policy and bargaining space may fail to identify cleavages and dimensions, decision-making analyses risk making false predictions, and – as de Gaulle's politics of the empty chair nicely exemplifies – the disregard of a country's vital interest expressed in a missing position could mislead our historical insight into a nation's power and sovereignty.

In our research, missing information on actors' positions could have these negative implications and affect our understanding of the deliberative process of constitution-building in the European Union (EU). Despite these negative effects, political scientists rarely discuss concepts dealing with missing data. In a broader examination of survey studies, King et al. (2001) found that, in 94% of the articles published between 1993 and 1997 in APSR, AJPS and BJPS, political scientists used listwise deletion of missing data, reducing their sample on the average by one-third. In this article, we want to draw attention to the more specific problem of missing positions and reveal how critical the decision either to delete or to impute missing positions is.

Outside of political science, generations of statisticians have developed unbiased, model-specific solutions designed to eliminate the loss of existing information. Under specific assumptions about the missing data-generating process, modern statistics go far beyond the conventional methods of missing value imputation, such as listwise deletion, pairwise deletion or mean imputation. Alternatives are based on maximum likelihood or multiple imputation using Bayesian statistics (Allison, 2002; Graham and Schafer, 2002). At the same time, the measuring of positions is becoming an increasingly common practice in political science (Laver and Hunt, 1992; Poole and Rosenthal, 1997; Budge et al., 2001; Laver, 2001). Scholars use actors' positions to make inferences about the nature and dimensionality of policy and bargaining spaces as well as to test prominent theories of political decision-making (Laver and Shepsle, 1994; Bueno de Mesquita and Stokman, 1994; König and Pöter, 2001; and Thomson and Stokman, 2006). Deleting dimensions or excluding actors because of missing values might significantly bias their analyses and findings.

In the following, we compare different methods for imputing missing

values within actors' positions. Our missing value problem concerns item non-response (estimators for positions on a few issues that are missing), which is generally solved by imputation. However, although research on missing values emphasizes the superiority of multiple imputation methods against case deletion, few studies have taken the problem of non-ignorable missing values into account – positions that are hidden for strategic purposes.¹ We attempt to consider this problem and discuss the advantages and drawbacks of (i) multiple imputation as implemented in AMELIA (King et al., 2001); (ii) a computationally easy but, in the context of spatial modelling, straightforward measure of indifference; and (iii) a conditional averaging algorithm, which might help scholars to reasonably fix the actors' positions in the policy space.

Criteria and goals of missing value imputation of positions

In any empirical analysis, we will not find all information stored in documents nor will interviewees answer all questions. In political science, one reason for non-documented information or refusal is that actors do not want to uncover their positions, which can be defined as actors' bargaining or ideal points (Bueno de Mesquita, 2004). For this reason, opinion polls and experts are often used to estimate the positions of the key political decision-makers. Apart from the problem of receiving reliable opinion estimators and finding expertise, some of the interviewees may still fail to answer a question on a particular topic. As a result, the data may contain various codes indicating the lack of an actor's position: 'don't know', 'refuse', 'unintelligible'. Once the study has been completed, all of these various answers have the same effect: they generate missing values on the positions of actors.

The question of whether missingness creates serious problems for political analysis easily prompts the answer 'yes, it does'. Most political science research and theories assume 'completeness', i.e. in terms of covering the whole policy and bargaining space or information, including the knowledge about all actors' positions. To impute missing values in a straightforward manner raises, however, methodological problems and analytical challenges. Note that we do not intend to recover, restore or predict missing values. The imputation of missing values is a means of achieving the overall goal of social research, namely 'to make valid and efficient inferences about a population of interest – not to estimate or predict, or recover missing observations nor to obtain the same results that we would have seen with complete data' (Graham and Schafer, 2002: 149). Our central goal is to impute missing values

in a way that minimizes distortions and bias of the respective analysis's inferential parameters (e.g. equilibria, coefficients).

In general, political scientists have a finite data set that consists of a representative sample of the population it was drawn from. These data usually contain a number of variables $V = \{v_1 \dots v_n\}$, and we seek to detect systematic and causal relationships between these variables (Bankhöfer, 1999). A necessary condition for generality, however, in making statistical inference above the finite number of outcomes is that the data set is representative of the population of outcomes. Compared with data used for statistics, data on actors' positions complicate the situation, although the goal of political analysis remains the same: most often the models are strategic and interactive, and the processes are dynamic.

Our DOSEI (Domestic Structures and European Integration) study attempts to collect information on the positions of the actors involved in the process of EU constitution-building. Knowing the actors' utility functions, their action sets and the structure of the political game, we might be able to make outcome predictions by solving the game for its equilibria. Thus, if we do not capture all relevant actors and their positions, the utilities and action sets should at least mirror the political game as played between the population of actors. Hence, in the event of missing positions, the goal is to keep the parameter estimates and the solutions to the game (e.g. core, yolk, finite set) undistorted as compared with the game played in the population of actors (or as compared with the population of games). But the best means of achieving this goal depends on the types of missing value in our DOSEI data set.

Types of missing value: Missingness

In current research on missing values, 'missingness' is regarded as a probabilistic phenomenon (Rubin, 1976). For any rectangular data set having the form: actors n times issues m , we can define variables that account for what is missing ($R_{nm} = 0$) and for what is observed ($R_{nm} = 1$). R describes a statistical distribution of missingness, even if we do not find a specific distribution. Because missingness may be related to the data, we can classify R according to the nature of that relationship between observed data and missing values (Graham and Schafer, 2002: 151). According to Rubin (1976), the main assumptions in the treatment of missing values are whether these values are missing completely at random, missing at random, or non-ignorable.

Political science studies mostly assume values to be missing completely at random (MCAR), which means that the probability for $R_{ym} = 0$ is independent from both the values on Y and the values on X (King et al., 2001).² Consider a

situation in which the positions on X ($X_1 \dots X_k$) are known for each actor, but the positions on Y ($Y_1 \dots Y_k$) are partly missing. MCAR is only given when the probability for missing values on Y depends neither on the other values of Y nor on the values of X (and, by independence, not on the values of other actors). In the very few political situations in which the missing values are completely at random, 'the set of individuals can be regarded as a simple random subsample from the original set of observations' (Allison, 2002: 3). Only if values are completely at random could one run the analysis on the set with complete information without having to worry about biased results.

Most of the MCAR critics claim that listwise deletion is highly problematic owing to the exclusion of cases (issues) with missing positions. The strong assumption holds only when the cases with complete information represent a random subsample within the original sample (Little and Rubin, 1987). King et al. (2001) find that 'the point estimate in the average political science article is about one standard error further away from the truth because of listwise deletion' (King et al., 2001: 52). The main virtue of listwise deletion is its simplicity, but, in our case, it would be likely to lead to false results. Since the deletion of a decisive issue or actor may create serious problems, it is preferable to draw on the weaker missing at random (MAR) assumption and to impute the missing values in Y using its systematic relation to X .

MAR assumes that the probability of missing data on a specific policy issue is independent of the positions of this issue once the effects of the remaining variables in the data set are taken into account. This means that the probability that $R_{ym} = 0$ may depend on X but not on Y . Although the deletion of cases is justifiable only under MCAR, the MAR assumption allows a researcher to use information about the relationship of Y to the other variables X in order to impute missing values. This relationship can be:

- the multivariate distribution of the variables,
- the distribution of the missing values across variables,
- the covariance and correlation between the variables,
- and, for data on positions, we can add the positional configuration (preference profiles) of the actors.

A whole battery of solutions exists for using the information from the other variables. Either multiple imputation (MI)³ or Maximum Likelihood (ML) imputations are recommended for linear analytical models (Allison, 2002; Graham and Schafer, 2002). ML has rarely been used in political science research because it is based on the principle of choosing 'as estimates those values that, if true, would maximize the probability of observing what in fact has been observed' (Allison, 2002: 13), given a specific model.

The third and possibly most realistic assumption is that non-ignorable missing data exist if the probability of missing values on a certain issue is dependent on the positions of this issue (MNAR). MNAR means that the probability that Y is missing depends on Y – an actor does not mention her position owing to other actors' positions. Here, some residual dependence between missingness and Y remains after accounting for X (Graham and Schafer, 2002: 151). MNAR requires a researcher to consider this selection bias explicitly within the analytical framework (Heckman, 1976; Winship and Mare, 1992; Agarwal, 2001). Dealing with data on actors' positions and strategies, we still might be suspicious about whether or not the data fulfil the MAR assumption or whether they suffer from non-ignorable missingness. We interviewed many experts from within and outside government, but some missing values still exist in the DOSEI data. To deal with this missing value problem, we will present the pros and cons of three prominent methods of imputation and then analyse these missing values in more detail.

Methods to impute missing positions

As for most political science research, deleting cases is no solution for our (and others') missing value problem. But how can we deal with missing positions? A more general recommendation would be to replace them either with another value or with the means (conditional or unconditional) of the other actors' positions. We believe that these solutions are also problematic because they dramatically reduce the variance of the variables with missing data by using the same value for all cases with missing data (Graham and Schafer, 2002: 160). Moreover, it is inappropriate to use mean imputation with variables measured at the nominal or ordinal level. They seriously distort inferences (tests and confidence intervals) by creating bias and overstating precision; thus, such imputations cannot be generally recommended (Little, 1992: 1231).

Sophisticated solutions to the missing value problem drawing on the MAR assumption usually argue that there is always uncertainty in missing data imputation, and by imputing only one value this uncertainty is artificially reduced to zero. As a result of eliminating the uncertainty from the model, the standard errors of the estimated coefficients are biased towards zero, making it easier to find significant relationships in the data. Next to ML estimation of missing values, multiple imputation is another way to address the uncertainty issue directly. It should be noted that, although the multiple imputation procedure assumes that the data are jointly multivariate normal, this assumption is very robust to departures from normality. Another advantage of this method is that the number of data sets imputed is relatively low:

'the relative efficiency of estimators with m as low as 5 or 10 is nearly the same as with $m = \infty$, unless missingness is exceptionally high' (King et al., 2001: 56). Certain requirements must be met for MI to have these desirable properties: the data must be missing at random, the model used to generate the imputed values must be correct, and the model used for the analysis must match the model used in the imputation (Rubin, 1976, 1987).

In general there are three facets of information in the data that can be utilized to reasonably impute the missing values in this context by considering (i) the entire data set (distribution, covariance etc. across cases and variables); (ii) the information in the single variables (issues); (iii) the information on the cases (actors). We will present three different methods for missing value imputation, each with a different focus on the observed data, and evaluate their usefulness within our context.

AMELIA

AMELIA offers a modified version of the expected maximization (EM) algorithm for multiple imputation (for details, see King et al., 2001). Multiple random imputation of missing values specifies the multivariate distribution of the variables in the data based on the existing data (see Rubin, 1987). The means and covariance are used to predict the regression coefficients (in the case of a linear regression model) for each of the variables. Using these regression coefficients, the missing values are imputed and the multivariate posterior distribution is specified. These steps are repeated until convergence is reached. The researcher can choose the number of complete data sets to be generated. All complete data sets are then combined according to Rubin's rule (1976) to calculate the statistical model of interest (see Allison, 2002: 29ff.; King et al., 2001).⁴

Compared with ML imputation, multiple imputation is 'probably less sensitive to the choice of the model because the model is only used to impute the missing data not to estimate the other parameters' (Allison, 2002: 32), but it nonetheless hinges on assumptions that might be critical for data applied to decision-making theories:

- all variables have normal distributions (although AMELIA appears to be fairly robust to deviations from normality; King et al., 2001),
- each variable can be represented as a linear function of all the other variables, together with a normal, homoscedastic error term (Allison, 2002).

Regarding non-ignorable missingness, AMELIA does not cover the likelihood function between the observed and the missing values on Y . Moreover, the

assumption of mutual linear representation between all variables is problematic if we consider different types of actor (such as domestic versus EU level) and different underlying policy dimensions. Although AMELIA offers solutions for the missing value problem in a straightforward and statistically precise mode of imputation, the problem of strategic hiding of positions remains and has to be addressed in the analytical model.

Indifference

As described by Hinich and Munger (1997), an efficient model of missing value imputation does not necessarily have to mirror the exact positional configuration of the complete political game. Instead it is sufficient that the imputed values for the actors do not distort the key solution concepts. This may imply that the reason for the missing position is the indifference of the specific actor with regard to the issue at stake. In other words, the solution concept of choice could be determined without knowing the exact position of every actor, but the analytical model requires that all actors have to be included. Indifference means that 'indifferent actors take a position such that the decision outcome agreed on by the other actors is allowed to pass; in other words, they interfere as little as possible' (Thomson and Stokman, 2006: 20). The underlying behavioural assumption is that 'actors who have no interest in a particular issue want to ensure that the outcomes of the discussions between the actors who do have an interest are reflected in the decision outcomes' (p. 20). Using indifferences for imputation also requires awareness of the underlying analytical model. For spatial models that will be applied to the DOSEI data, it might be reasonable to impute the middle (indifference) position between the status quo (or reference point) and the draft constitution. As far as possible win sets to the status quo are concerned, this imputation method ensures that no possible solution between the analytical fix points is diminished or reduced.

Although indifference is straightforward for spatial analysis, its drawbacks are obvious:

- the model hinges upon the assumption that missing positions are due to the actors' indifference between two analytically derived points in the policy space,
- in order to make use of this idea in the most efficient way, it should be implemented on a model-specific basis. This can pose challenges to cross-model comparison.

In the light of a non-ignorable missing value problem, indifference is possibly

problematic because hiding positions suggests a more extreme position. It is therefore questionable that indifference is the primary explanation for missingness in the DOSEI data if the MNAR assumption is met.

Likelihood to deviate from the mean (LDM)

Another approach to evaluating the goodness of an imputation model applied to a positional data set compares the actors' configuration in an n -dimensional policy space in the data set before and after imputation. Assuming that each variable on positions is of equal weight, we try to minimize the following equation:

$$\frac{MIN}{i=1} \sum_{j=2}^{n-1} \sum_{j=2}^n |E_{ij}^t - E_{ij}^{t-1}| \quad (1)$$

where E is the mean Euclidian distance between actor i and actor j across m variables k before ($t-1$) and after the imputation (t):

$$E_{ij} = \sqrt{\sum_{k=1}^m (v_{kj} - v_{ki})^2}, \text{ where } j \neq i. \quad (2)$$

Trying to solve the minimization problem, we are left with a system of equations for which no unique solution exists. Assuming completeness and the equal weight of all issues, we first need to order the issues in accordance with an underlying policy dimension, such as more vs. less integration, left vs. right (Zimmer et al., 2005). Ordinal data are then transformed on a pseudo-metric scale via marginal normalization⁵ and standardized on an interval $\{0;1\}$ with $\mu = .5$. One obvious solution to equation (1) would be to use the mean position an actor holds across all issues referring to the same underlying policy dimension. Owing to the marginal normalization applied to the data before, this corresponds to the ML estimator for the joint cumulative probability function of actor i and variable x . However, because we are using pseudo-metric data only, there may be no equivalent to the mean on the underlying ordinal scale. Therefore we take the value of the underlying ordinal scale that is most likely (in terms of cumulative probability under the p.d.f). To fit the imputation model with the configuration of the completed data set, we repeat the whole process (using the completed data set) until convergence.⁶

This algorithm may be best understood as a combination of conditional averaging across the single cases (Graham and Schafer, 2002: 157f.) and a case-wise ML estimation. It suffers from two major drawbacks:

- the assumption that all issues are of equal importance for the underlying policy dimension, and
- the ordering of all issues in accordance with an underlying policy dimension.

The LDM concept might partially solve a censoring problem in the DOSEI study. Based on the knowledge about the actors' general stance towards the underlying policy dimension, LDM imputes a reasonable value for the sensitive issues that one cannot observe. In contrast to AMELIA, it does not require mutual linear relationships between the variables and is not distorted by the different functional roles that actors may perform in the decision-making process. If we want to dismiss the rigorous assumption that all issues are of equal weight, we could combine LDM with a latent trait model found in Item Response Theory (Congdon, 2001: 223 ff.).

Missingness in data on positions: DOSEI data

What kind of information did we gather for measuring positions in the DOSEI study? If missing values occur, which assumptions (MCAR, MAR, MNAR) can we make about the (missing) data-generating process, and which method can we use to impute missing positions? To answer these questions, we turn our attention to the research design and sampling techniques, which are related to the missing value problem. In the DOSEI study, we use data gathered by the following instruments:

- text/content analysis (*inter alia*, Laver and Budge, 1992; Klingemann et al., 1994; Budge et al., 2001; Gabel and Hix, 2002; Laver et al., 2003; Pennings et al., 2004),
- expert interviews (*inter alia*, Castles and Mair, 1984; Laver and Hunt, 1992; Bueno de Mesquita and Stokman, 1994; Huber and Inglehart, 1995; Ray, 1997; Thomson and Stokman, 2006), where sampling often depends on how much leeway the researcher is willing to grant the expert,
- survey data using opinion polls (*inter alia*, Gerber, 1996; König and Hug, 2000),
- statistical analyses of roll-call votes (e.g. Poole and Rosenthal, 1997; Cox and Poole, 2001; Mattila and Lane, 2001; Mattila, 2004), and
- delegates and parliamentary amendments (Tsebelis and Money, 1995).

Although the number of missing positions may vary across the data gathered by these instruments, they share one common problem: a conscious

and subjective choice of actors and/or policy issues by either the researcher, the public or an expert remains vulnerable to selection bias. The limited availability of information on the positions brings about a risk of non-ignorable missing values. Our strategy of interviewing more than one expert per country is one reason for the small percentage of missing values observed in the aggregated DOSEI data set. We asked the experts to indicate only the positions of actors deviating from the national position (NP), which allows us to replace non-indicated values for domestic actors with the respective NP. The rectangular data set contains 136 actors with 65 questions on the relevant issues, and thus 8840 cells, of which 389 are not filled. From these 389 cells with no position, 252 could be filled with the NP indicated by the experts. However, 137 (or roughly 1.6%) of the data entries still remain as missing values. As Table 1 shows, their distribution varies – with very few exceptions – at the level of entire countries but not at the level of actors.

Explaining the occurrence of missing positions

According to social psychologists (Tourangeau et al., 2000; Tourangeau and Smith, 1996), an interviewee has to (1) understand the question, (2) recall the relevant information, (3) map this information onto the possible answers, and (4) be willing to edit the correct answer. In terms of missing positions, whereas the first three steps correspond to an interviewee's lack of competence, the fourth step concerns unwillingness or inability to reveal the correct answer. We therefore propose to check the appropriate assumption on the probability of missing positions in Y in the DOSEI data set by distinguishing between two reasons for missing values:

- (a) the expert did not know the position of actor i on issue Y for any reason other than (b) (*lack of competence*),

Table 1 Missing policy positions in the DOSEI data set (total $N = 8840$)

	<i>AP observed</i>		<i>AP missing</i>	
	<i>N</i>	%	<i>N</i>	%
NP observed	8451	95.6	252	2.85
NP missing	8	0.09	129	1.46

Note: 27 official national positions [NP] (member states, Commission, EP)
+ 109 domestic actors' positions [AP] = 136 actors \times 65 issues = 8840 cases.

- (b) either the expert hid extreme positions on sensitive issues or the actor was able to disguise the position for *strategic purposes*.

In the event of the first situation, we are close to fulfilling the MCAR or at least the MAR assumption because the limited knowledge of our experts on Y may reveal some covariance structure with X . However, if the second situation is true, we can rely only on MAR or we may even have to deal with MNAR. In this case, it is plausible to suppose that an actor strategically hides her position on Y because of the positions held by other actors on Y , so that the probability of $R = 0$ depends on the values of Y . Whether we should assume MNAR to be more appropriate for defining how the missing values on Y might look is, however, an 'extremely difficult task' to solve (Allison, 2002: 71). But, if actors hide positions for strategic purposes, the question arises of how this aspect of the missing data-generating process should be considered in our analysis – even if the percentage of missing values is comparatively small.

A reasonable starting point for thinking about the occurrence of missing positions might be to ask under what conditions it is advantageous for an actor to hide her national position on a specific issue. In our study, we can assume that all countries must agree on the constitutional text, which might imply a bargaining advantage for actors biased towards the status quo. They can expect to gain concessions, whereas actors more supportive of policy change have to make offers in order to find agreement for changing the status quo. In these circumstances, if an actor is located close to or at the status quo of an issue, it does not make sense for her to hide her position. This actor will gain concessions only when the other actors share a more supportive position and know about her status quo bias. Hiding a status quo position accordingly risks only that the other actors might wrongly believe that she is in favour of the more supportive position. Similarly, if this actor shared the views of the other actors, why should she risk that they might possibly draw another (wrong) conclusion and change the outcome towards a more extreme position? There seems to remain only one condition under which an actor will reasonably hide her position, namely when the other actors need to be led to believe that she has a status quo bias on the draft text and are thus willing to offer her concessions, when she in fact has a more supportive position on the specific issue. Only in this case must this actor fear gaining fewer concessions and therefore have an incentive to conceal her position. This suggests a selection bias of missing values for status quo located actors far away from other actors' (supportive) positions.

Examining missing positions with selection bias

The usual suggestion for linear modelling is to include a possible selection bias from the underlying structure in the analytical model (e.g. Heckman, Tobit; see Winship and Mare, 1992). One way to uncover the underlying structure of missingness is to estimate a probit regression on whether or not a case had any missing values, followed by a Poisson regression on the number of missing values encountered. Furthermore, one can separate both effects by estimating a Heckman full ML selection model, which includes three types of variables for our analysis:

- (1) An actor's mean distance to the status quo (MESQ) as a proxy for her status quo bias, and an actor's mean distance to the mean of other national positions (MENP) as a proxy for the distance to the bargaining outcome at the EU level. A large distance to NP and a small distance to SQ would confirm our expectation that actors with strategic views on constitution-building explain the occurrence of missing positions.
- (2) Variables on the actors' position in relation to the domestic political debate: a measure for heterogeneity operationalized as an actor's mean distance to the other actors in her country (AVDIST), an actor's distance to the national position (DISTNP), an actor's average greatest distance to other actors of the same country (DISTMAX), and the average percentage of issues where at least one domestic actor's position is identical to the status quo (PERSQ). These variables refer to the domestic situation and indicate whether the relative location of an actor, with respect to her preference of the status quo, her distance from other domestic actors, the national position and the opposite actors, explains the missing positions.
- (3) Finally, we add the number of interviewees per actor found in a country as a technical control variable (NUMINT).

We adjust all variables for effects resulting from different scales.

Using the DOSEI data set we first estimate a probit regression on a dummy variable indicating whether or not we encountered at least one missing value for a particular case. According to Table 2, MESQ and MENP reveal large and significant effects on the likelihood of encountering any missing values: the larger the distance to the approximate bargaining outcome (estimated as the mean of all national positions) and the smaller the distance to the status quo, the more likely it is we will encounter at least one missing position. We also find that this likelihood increases with the contestation of the national position (NP) within a country (DISTNP).

Next we estimate a Poisson regression on the number of missing values

Table 2 Probit model, Poisson regression and Heckman Selection Model on missing policy on the occurrence and number of missing policy positions

	<i>Probit</i>		<i>Poisson</i>		<i>Heckman</i>	
	<i>Coef.</i>	<i>S.E.</i>	<i>Coef.</i>	<i>S.E.</i>	<i>Coef.</i>	<i>S.E.</i>
	<i>N</i> = 136		<i>N</i> = 136		(uncens. obs. = 34)	
Distance to NP (DISTNP)	3.29**	-1.09	6.89***	-2.35	7.24***	-1.15
Maximum distance to other domestic actor (DISTMAX)	0.33	-1.23	-3.99***	-0.91	-7.55***	-1.42
Average distance to other domestic actor (AVDIST)	3.38	-3.70	-2.18*	-1.31	-16.21***	-4.48
Number of interviewees (NUMINT)	0.13	-0.21	-0.34**	-0.15	-0.31	-0.25
Percentage of issues with IP = SQ (PERSQ)	-0.86	-1.94	11.39**	-2.85	-6.62***	-1.63
Mean distance to the SQ (MESQ)	-2.11***	-0.67	-2.46***	-0.44		
Mean distance to mean of all NPs (MENP)	2.37**	-1.36	8.04***	-2.47		
Constant	0.22	-2.06	19.89***	-3.99	10.61***	-1.04
<i>Selection stage</i>					(cens. obs. = 102)	
Mean distance to the SQ (MESQ)					-0.93***	-0.27
Mean distance to mean of all NPs (MENP)					3.86***	-0.95
Constant					-1.60**	-0.80
athrho					-1.45	-0.65
lnsigma					0.53***	-0.21
rho					-0.89	-0.13
sigma					1.70	-0.35
lambda					-1.52	-0.50
LR test of independence of equations (rho = 0): Prob > chi ² =						.001
Wald/Chi/LR test: Prob > chi ² =	.000		.000		.000	
Pseudo- <i>R</i> ²			.692		.315	

*significant at 10% level; **significant at 5% level; ***significant at 1% level.

encountered. This analysis confirms the effects as already revealed by the probit regression plus additional effects with respect to the domestic constellation of positions: the larger the mean distance of an actor to the other domestic actors (AVDIST) and the larger the average largest distance to any other domestic actor (DISTMAX), the fewer missing values we can expect. The results confirm our expectation and reveal that two different effects are responsible for missing values: the first is the strategic hiding of positions where these actors have a status quo bias and are far from the other actors' mean position. This effect appears to be especially related to the positions held at the EU level. The second effect concerns the prominence of a position and thus the ease of identification by the interviewee.

In order to separate the effects on the likelihood and the number of missing positions, we finally run a full ML Heckman model. The two variables related to the policy space are included in the selection stage and we insert the other variables into the second stage. The coefficients reveal the same algebraic sign as in the Poisson regression, but the coefficients are slightly larger. The percentage of issues where at least one of the actors (in the same country) has a position identical to the status quo has a negative impact on the number of missing values. One possible reason for this relationship is that it is easier to identify positions during the national debate if one of the relevant actors is on the status quo. The value for ρ ($= .89$) indicates the expected selection bias between these two effects on this level of data aggregation, although we cannot completely avoid correlation in the residuals of both steps.⁷

The empirical analyses show that the different actors' locations in the policy space determine the likelihood of encountering any missing values in the first place. Actors strategically hide their position when they can expect to receive more concessions. The number of missing positions is additionally explained by the ease of identification or the prominence of an actor's position compared with that of other domestic actors. This suggests that the remaining missing values in the DOSEI data set are MAR, if not MNAR. But how do the three methods cope with this problem? Do they vary in their imputation results and, if so, to what extent do they provide us with a plausible imputation?

Variance and goodness of imputation

To answer these questions we first assess whether and to what extent the different modes impute different values, focusing on the differences between AMELIA, LDM and indifference. Table 3 displays the number of positions imputed differently and the average difference on a pseudo-metric scale by

Table 3 Variance in imputation: Number of differently imputed values and average difference on pseudo-metric scale

	<i>LDM</i>	<i>AMELIA</i>	<i>Indifference</i>
LDM	0	103	68
Amelia	93.6%	0	114
Indifference	42.3%	96.1%	0

Note: *N* = 129.

each combination of the three methods. As the underlying ordinal scale of the 65 issues has 2.8 real values on average, it is remarkable that AMELIA imputes 114 and 103, respectively, of the 129 missing values differently from the imputations suggested by indifference and LDM. LDM and indifference, by contrast, produce only 68 different values.

The three imputation methods produce quite different solutions and it remains difficult to find an absolute criterion for the goodness of the imputations in a methodologically rigorous manner. We propose to rely on our theoretical and empirical analyses of the underlying mechanism suggesting that actors biased towards the status quo and far from other actors may have an increased motivation strategically to conceal their more supportive positions. Thus, when we presume that the expected outcome may be approximated by the mean position of all actors, we are able to distinguish the interval between status quo and expected outcome (status quo interval) from the interval between expected outcome and draft proposal (draft interval) in order to identify a more supportive hidden position.

For these conditions, Table 4 lists a ‘hit rate’ of imputed values (IV), where $IV > \text{expected outcome}$. Compared with indifference, which systematically imputes more values in the status quo interval (because the expected outcome is closer to the draft than to the status quo), AMELIA imputes slightly more positions in the draft interval. Finally, LDM reveals the highest ‘hit rate’ of a more supportive position on the draft. The reason is that, across all

Table 4 Hit rate ($|\text{Imputed Value} - \text{Draft}| < |\text{Expected Outcome} - \text{Draft}|$)

<i>Scale</i>	<i>LDM</i>	<i>AMELIA</i>	<i>Indifference</i>
Ordinal scale	82.7%	64.6%	35.4%
Pseudo-metric scale	78.0%	59.1%	55.9%

dimensions, even the actors relatively biased towards the status quo are still close to the draft position, which frequently lies in the 'draft interval'. Using an ordinal scale, LDM imputes about 74% of all cases of the draft position.

Conclusion

Missing or incomplete data on actors' positions can cause significant problems in political analysis. Although this problem has rarely received attention in political science research, the growing number of political analyses increases the relevance of the data problem of missing policy positions. Political analysis uses actors' positions as the basic variable for the study of political phenomena, and theories and tools of analysis are not designed to ignore missing values. In our view, when we use data such as the DOSEI data on constitution-building for explorative, comparative and theory-testing purposes, 'application-specific approaches' are often 'worth the trouble' in handling missing values (King et al., 2001: 58). In our view, the problem of non-ignorable missing positions (MNAR) is non-ignorable.

In political science and our research on constitution-building, we can hardly assume non-strategic behaviour, but the imputation for MNAR requires a model-specific solution. If the model is linear, Heckman and Tobit are possible solutions to this problem. But, if the model is not linear, no standard tool is available, although the advances in Item Response Theory suggest promising concepts (Winship and Mare, 1992; Congdon, 2001: 233 ff.). Instead of using stronger assumptions, such as MCAR and MAR, we propose examining the likelihood of the strategic concealment of sensitive positions in the data set, even though we encounter a very small percentage of missing positions in the DOSEI data.

Our findings suggest that we can neither use MCAR and delete the few cases with missing positions, nor completely reject MNAR. Missing positions occur when actors biased towards the status quo seem to have extreme positions vis-à-vis the other actors at the EU level, and the number of missing positions is explained by the identification problem of the interviewee. Although none of the existing methods for imputing missing values is perfectly suited to MNAR, LDM seems to come close to resolving our missing value problem in the DOSEI data. Compared with the prominent AMELIA tool, LDM considers the relative positions of the individual actor. Indifference is based on an opposite theoretical assumption about the underlying mechanisms leading to missing values.

Our empirical evaluation also shows that the three concepts produce quite different results. In other words, researchers should carefully consider

which method of imputation they wish to apply for imputing positions. Even with a comparatively low percentage of missing positions, the method of imputation can lead to significantly different findings because the resulting imputed values may vary greatly. We recommend (a) checking carefully the underlying causes of missingness and examining whether MCAR, MAR or MNAR may hold; (b) considering whether political analysis needs a complete data set by taking into account the theoretical insights on the policy space; (c) deciding which concept of missing value imputation solves the missing value problem best; (d) bearing in mind the missing value problem when interpreting the findings of the policy analysis.

Notes

The authors are grateful for helpful comments by Simon Hug and financial support of DOSEI by the European Commission under grant No. HPSE-CT-2002-00117.

- 1 The literature dealing with non-ignorable missing values generally recommends using selection models (e.g. Heckman, 1985; Winship and Mare, 1992) and pattern mixture models (Little, 1993). Most of it deals, however, with unit non-response in either longitudinal studies (e.g. Molenbergh et al., 1998; Glynn et al., 1986) or bias introduced by selection on Y (Hug, 2003).
- 2 If we separate our complete data Y_{com} into a set containing missing values (Y_{mis}) and a set with non-missing values (Y_{obs}), missing completely at random means that $P(R | Y_{com}) = P(R)$, and missing at random occurs when $P(R | Y_{com}) = P(R | Y_{obs})$.
- 3 In particular, MICE (Multivariate Imputation by Chained Equations) for R is another freely available and useful MI package (<http://web.inter.nl.net/users/S.van.Buuren/mi/html/mice.htm>). For an overview of other, primarily commercial, programs for MI see Allison (2002: 34).
- 4 The algorithms leading to convergence can be divided into data-augmentation algorithms – especially IP (Imputation-Posterior) algorithms – and EM algorithms. Whereas IP takes random draws from the entire posterior in order to calculate the subsequent regression coefficients, the EM uses the deterministic means of the posterior. Data augmentation is slightly more accurate in theoretical terms, but it is computationally much more involved and suffers from convergence problems. AMELIA uses EMis (Expected Maximization sampling importance/resampling) as a default algorithm. This algorithm as employed by King et al. (2001) steps back in estimation uncertainty by sampling from the asymptotic approximation of the entire posterior and controlling for what the authors call the ‘importance ratio’ – the ratio of the actual posterior to the asymptotic approximation. EMis seeks to combine computational efficiency and accuracy. AMELIA, including documentation, can be downloaded from URL: <http://G.King.harvard.edu/>.
- 5 Marginal normalization divides the space under the standard normal

distribution into k equal parts in accordance with the relative frequency of each class k (for details, see Hartung, 2004). It is equivalent to a probit link function in Item Response Theory (e.g. Johnson and Albert, 2001).

- 6 If we wished to keep the parameters of the p.f. stable, we could use the posterior p.d.f.s. In this case the algorithm would be close to EM, but based on a different model.
- 7 We repeated the same analyses using the raw data. The algebraic signs of the coefficients remained the same, but the pseudo- R decreased significantly. In addition, it was impossible to separate the effects in a Heckman selection model ($\rho = 1$). We conclude that this difference is due to a higher percentage of missingness at random.

References

- Agarwal, Sameer (2001) 'Learning from Incomplete Data', URL (consulted September 2004): <http://www.cs.ucsd.edu/users/elkan/254spring01/sagarwalrep.pdf>.
- Allison, Paul D. (2002) *Missing Data: Quantitative Applications in the Social Sciences*. London: Sage Publications.
- Bankhöfer, Udo (1999) *Unvollständige Daten- und Distanzmatrizen in der Multivariaten Datenanalyse*. Cologne: Verlag Josef Eul.
- Budge, Ian, Hans-Dieter Klingemann, Andrea Volkens, Judith Bara and Eric Tanenbaum (2001) *Mapping Policy Preferences: Estimates for Parties, Electors and Governments 1945–1998*. Oxford: Oxford University Press.
- Bueno de Mesquita, Bruce (2004) 'Decision-Making Models, Rigor and New Puzzles', *European Union Politics* 5(1): 125–38.
- Bueno de Mesquita, Bruce and Frans Stokman (eds) (1994) *European Community Decision Making: Models, Applications and Comparison*. New Haven, CT/London: Yale University Press.
- Castles, Francis G. and Peter Mair (1984) 'Left–Right Political Scales: Some "Expert" Judgments', *European Journal of Political Research* 12(1): 73–88.
- Congdon, Peter (2001) *Bayesian Statistical Modelling*. New York: Wiley.
- Cox, Gary W. and Keith T. Poole (2001) 'On Measuring Partisanship in Roll Call Voting: The U.S. House of Representatives. 1877–1999', Mimeo, University of California at San Diego.
- Gabel, Matthew and Simon Hix (2002) 'Defining the EU Political Space: An Empirical Study of the European Election Manifestos, 1979–1999', *Comparative Political Studies* 35(8): 934–64.
- Gerber, Elisabeth R (1996) 'Legislative Response to the Threat of the Popular Initiative', *American Journal of Political Science* 40(1): 99–128.
- Glynn, R.J., N.M. Laird and D.B. Rubin (1986) 'Selection Modeling versus Mixture Modeling with Nonignoreable Nonresponse', in H. Wainer (ed.) *Drawing Inference from Self-Selected Samples*, pp. 115–51. New York: Springer-Verlag.
- Graham, Joseph L. and John W. Schafer (2002) 'Missing Data: Our View of the State of the Art', *Psychological Methods* 7(2): 147–77.
- Hartung, Joachim (2004) *Grundkurs Statistik*, 12th edn. Munich: Oldenbourg.

- Heckman, James (1976) 'The Common Structure of Statistical Models of Truncation, Sample Selection and Limited Dependent Variables, and a Simple Estimator for Such Models', *Annals of Economic and Social Measurement* 5: 475–92.
- Heckman, James (1985) 'Selection Bias and Self-selection', in Peter Newman (ed.) *The New Palgrave: A Dictionary of Economics*, pp. 287–96. New York: Macmillan Press.
- Hinich, Melvin J. and Michael C. Munger (1997) *Analytical Politics*. Cambridge: Cambridge University Press.
- Huber, John D. and Ronald Inglehart (1995) 'Expert Interpretations of Party Space and Party Locations in 42 Societies', *Party Politics* 1: 73–111.
- Hug, Simon (2003) 'Selection Bias in Comparative Research: The Case of Incomplete Data Sets', *Political Analysis* 11: 255–74.
- Johnson, Valen E. and James H. Albert (2001) *Ordinal Data Modeling*. New York: Springer.
- King, Gary, James Honaker, Anne Joseph and Kenneth Scheve (2001) 'Analyzing Incomplete Political Science Data: An Alternative Algorithm for Multiple Imputation', *American Political Science Review* 95(1): 49–69.
- Klingemann, Hans-Dieter, Roland Hofferbert and Ian Budge (1994) *Parties, Policies and Democracy*. Boulder, CO: Westview Press.
- König, Thomas and Simon Hug (2000) 'Ratifying Maastricht: Parliamentary Votes on International Treaties and Theoretical Solution Concepts', *Journal for European Politics* 1: 93–124.
- König, Thomas and Mirja Pöter (2001) 'Examining the EU Legislative Process: The Relative Importance of Agenda and Veto Power', *Journal for European Politics* 3: 329–51.
- Laver, Michael (2001) *Estimating the Policy Position of Political Actors*. New York/London: Routledge.
- Laver, Michael and Ian Budge (eds) (1992) *Party Policy and Government Coalitions*. London: Macmillan.
- Laver, Michael and William B. Hunt (1992) *Policy and Party Competition*. New York/London: Routledge.
- Laver, Michael and Kenneth A. Shepsle (1994) *Cabinet Ministers and Parliamentary Government*. Cambridge: Cambridge University Press.
- Laver, Michael, Kenneth Benoit and John Garry (2003) 'Extracting Policy Positions from Political Texts Using Words as Data', *American Political Science Review* 97(2): 83–112.
- Little, Rodrick J.A. (1992) 'Regression with Missing X's: A Review', *Journal of the American Statistical Association* 87: 1227–37.
- Little, Rodrick J.A. (1993) 'Pattern-Mixture Models for Multivariate Incomplete Data', *Journal of the American Statistical Association* 88: 125–34.
- Little, Rodrick J.A. and Donald B. Rubin (1987) *Statistical Analysis with Missing Data*. New York: John Wiley.
- Mattila, Mikko (2004) 'Contested Decisions: Empirical Analysis of Voting in the EU Council of Ministers', *European Journal of Political Research* 43:29–50.
- Mattila, Mikko and Jan-Erik Lane (2001) 'Why Unanimity in the Council? A Roll Call Analysis of Council Voting', *European Union Politics* 2(1): 31–52.
- Molenbergh, G., B. Michelies, M.G. Kenward and P.J. Diggle (1998) 'Monotone

- Missing Data and Pattern-Mixture Models', *Statistica Neerlandica* 52(2): 153–61.
- Pennings, Paul, Madeleine Hosli and Christine Arnold (2004) 'The EU Constitution and Positions on Governance', Paper presented at the ECPR Joint Sessions of Workshops, Uppsala, April 2004.
- Poole, Keith T. and Howard Rosenthal (1997) *Congress: A Political-Economic History of Roll Call Voting*. New York: Oxford University Press.
- Ray, Leonard (1997) 'Measuring Party Orientations towards European Integration', Mimeo, Department of Political Science, University of North Carolina.
- Rubin, Donald B. (1976) 'Inference and Missing Data', *Biometrika* 63: 581–92.
- Rubin, Donald B. (1987) *Multiple Imputation for Nonresponse in Surveys*. New York: Wiley.
- Thomson, Robert and Frans N. Stokman (2006) 'Research Design', in R. Thomson, F.N. Stokman, C.A. Achen and T. König (eds) *Decision Making in the EU*. Cambridge: Cambridge University Press, forthcoming.
- Tourangeau, Roger, Lance J. Rips and Kenneth Rasinski (2000) *The Psychology of Survey Response*. New York/Cambridge: Cambridge University Press.
- Tourangeau, Roger and T.W. Smith (1996) 'Asking Sensitive Questions: The Impact of Data Collection Mode, Question Format, and Question Context', *Public Opinion Quarterly* 60: 275–304.
- Tsebelis, George and Jeanette Money (1995) 'Bicameral Negotiations: The Navette System in France', *British Journal of Political Science* 25(1): 101–29.
- Winship, Christopher and Robert D. Mare (1992) 'Models for Sample Selection Bias', *Annual Review of Sociology* 18: 327–50.
- Zimmer, Christina, Gerald Schneider and Michael Dobbins (2005) 'The Contested Council: The Conflict Dimensions of an Intergovernmental Institution', URL (consulted March 2005): <http://www.jhubc.it/ecpr-bologna/docs/148.pdf>.

About the authors

Thomas König is Professor of Political Science, German University of Administrative Sciences, Freiherr-vom-Stein-Straße 2, 67346 Speyer, Germany.

Fax: +49 6232 654 127

E-mail: tkoenig@dhv-speyer.de

Daniel Finke is a PhD student at the German University of Administrative Sciences and a research fellow at the Research Institute for Public Administration, Freiherr-vom-Stein-Straße 2, 67346 Speyer, Germany.

Fax: +49 6232 654 127

E-mail: finke@foev-speyer.de

Stephanie Daimer is a PhD student at the German University of Administrative Sciences and a research fellow at the Research Institute for Public Administration, Freiherr-vom-Stein-Straße 2, 67346 Speyer, Germany.

Fax: +49 6232 654 127

E-mail: daimer@foev-speyer.de
